

## STATS315A PROBLEM SET 2

Due Date: Wednesday, March 6, 11:59pm

**Question 2.1** (Resampling methods and normal inference, 40 points): In this problem, we compare three parameter inference methods: (i) model-based inference, which assumes the model is true, (ii) the “sandwich estimator” we derived in class, and (iii) the bootstrap resampling estimator of variance. We will do this both for linear and (binary) logistic regression, repeating the following experimental protocol many times and providing summary results. We first describe the protocol for the Gaussian linear model case; we then describe modifications for the other cases.

i. For the data model

$$y_i = x_i^T \beta^* + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathbf{N}(0, \sigma^2),$$

generate a sample of size  $n$  (to be specified) in dimension  $d = 10$ , where  $x_i \stackrel{\text{iid}}{\sim} \mathbf{N}(0, I_d)$  and  $\sigma^2 = 1$ , and draw  $\beta^* \sim \text{Uni}(\mathbb{S}^{d-1})$ , the sphere in  $\mathbb{R}^d$ .

ii. Compute  $\hat{\beta}_n$  minimizing

$$L_n(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \beta)^2.$$

iii. Consider the following three covariance estimates:

$$\Sigma_n := \hat{\sigma}^2 (X^T X)^{-1}, \quad \hat{\sigma}^2 = \frac{1}{n-d} \sum_{i=1}^n (y_i - x_i^T \hat{\beta}_n)^2 \quad (\text{FISHER})$$

$$\Sigma_n := (X^T X)^{-1} \left( \sum_{i=1}^n (y_i - x_i^T \hat{\beta}_n)^2 x_i x_i^T \right) (X^T X)^{-1} \quad (\text{SANDWICH})$$

$$\Sigma_n := \frac{1}{B} \sum_{b=1}^B (\hat{\beta}_n^b - \hat{\beta}_n)(\hat{\beta}_n^b - \hat{\beta}_n)^T \quad (\text{BOOTSTRAP})$$

where  $\hat{\beta}_n^b$  is a bootstrap resampled least-squares estimate, and  $B = 200$  is the number of bootstrap replicates.

The first covariance is the classical Fisher information matrix, the second the covariance as per the standard asymptotic theory we have developed, and the third that of the bootstrap. We have  $\hat{\beta}_n \sim \mathbf{N}(\beta^*, \Sigma_n)$  for each of these (so long as the data model remains true!), and therefore

$$\mathcal{C}_n := \left\{ \beta \in \mathbb{R}^d \mid (\beta - \hat{\beta}_n)^T \Sigma_n^{-1} (\beta - \hat{\beta}_n) \leq \chi_{d,1-\alpha}^2 \right\}$$

as an asymptotically valid  $1 - \alpha$  confidence set, that is,  $\mathbb{P}(\beta^* \in \mathcal{C}_n) \rightarrow 1 - \alpha$ , where  $\chi_{d,1-\alpha}^2$  is the  $1 - \alpha$  quantile of a  $\chi^2$  random variable with  $d$  degrees of freedom. Use  $\alpha = .1$  for the remainder.

(a) Repeat the experiment **i-iii** for  $n = 50, 100, 200, 400$  for  $T = 200$  times, and track the fraction of times that  $\beta^* \in \mathcal{C}_n$  for each of the covariances (**FISHER**), (**SANDWICH**), and (**BOOTSTRAP**). For each covariance approximation, plot your coverage against sample size  $n$ .

- (b) Repeat part (a) except for logistic regression. Thus, make the following changes to the procedure i–iii. Instead of the linear regression model, generate data from the logistic regression model

$$\mathbb{P}_\beta(Y = y \mid X = x) = \frac{e^{y\beta^T x}}{1 + e^{x^T \beta}}, \quad y \in \{0, 1\},$$

where  $x_i, \beta^*$  are generated identically. In part ii, choose  $\hat{\beta}_n$  to minimize the negative log likelihood, that is, for  $\ell(\beta, x, y) = -\log p_\beta(y \mid x) = \log(1 + e^{x^T \beta}) - yx^T \beta$ , let  $L_n(\beta) = \frac{1}{n} \sum_{i=1}^n \ell(\beta, x_i, y_i)$ . In part iii, replace the Fisher information and Sandwich covariances with their counterparts

$$\begin{aligned} \Sigma_n &:= (X^T W X)^{-1}, \quad W = \text{diag}([\hat{p}_i(1 - \hat{p}_i)]_{i=1}^n) \\ \Sigma_n &:= \frac{1}{n} \nabla^2 L_n(\hat{\beta}_n)^{-1} \widehat{\text{Cov}}(\nabla \ell(\hat{\beta}_n(X, Y))) \nabla^2 L_n(\hat{\beta}_n)^{-1}, \end{aligned}$$

respectively (the bootstrap covariance does not change).

- (c) Repeat part (a) with faulty linear modeling assumptions, so that we evaluate coverage of  $\hat{\beta}_n$  and  $\mathcal{C}_n$  for the best linear predictor in mean-squared error,  $\beta^{\text{mse}} = \text{argmin}_\beta \mathbb{E}[(y - x^T \beta)^2]$ . To do so, replace the model in part i with

$$y_i = x_i^T \beta^* + (x_i^T \theta^*)^2 + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathbf{N}(0, \sigma^2),$$

where  $\theta^* \sim \text{Uni}(\mathbb{S}^{d-1})$  as well. Note that

$$\begin{aligned} \mathbb{E}[(y - x^T \beta)^2] &= \mathbb{E}[(\varepsilon + x^T(\beta - \beta^*) + (x^T \theta^*)^2)^2] \\ &= \sigma^2 + \|\beta - \beta^*\|_2^2 + \mathbb{E}[(x^T \theta^*)^4], \end{aligned}$$

because  $\mathbb{E}[v^T x (u^T x)^2] = 0$  for any vectors  $u, v$ ,<sup>1</sup> so  $\beta^{\text{mse}} = \text{argmin}_\beta \mathbb{E}[(y - x^T \beta)^2] = \beta^*$  as well.

**Question 2.2** (Asymptotics of causal inference, 15 points): As in Question 1.6, consider the potential outcomes framework for a real-valued response  $Y$  with randomized treatment assignments  $W \in \{0, 1\}$ , so that  $(Y(0), Y(1)) \perp W$ . Let the population mean-square-error estimates be

$$(\tau^{\text{mse}}, \alpha^{\text{mse}}, \beta^{\text{mse}}) = \text{argmin}_{\tau, \alpha, \beta} \mathbb{E}[(Y - \alpha - X^T \beta - \tau W)^2],$$

so that  $\alpha^{\text{mse}} \in \mathbb{R}$  is an intercept,  $\beta^{\text{mse}} \in \mathbb{R}^p$ , and  $\tau^{\text{mse}}$  is the coefficient of  $W$  in the model  $Y_i = \alpha + X_i^T \beta + \tau W_i + \varepsilon_i$ . Given a sample of size  $n$ , where each individual is chosen to be in treatment ( $W = 1$ ) or control ( $W = 0$ ) with equal probability  $\frac{1}{2}$ , let  $\hat{\tau}_n, \hat{\alpha}_n, \hat{\beta}_n$  be the empirical squared error minimizers. Give the limiting distribution of  $\hat{\tau}_n$ . That is, give the value of the asymptotic variance  $\sigma^2(\tau)$  in the limiting normal

$$\sqrt{n}(\hat{\tau}_n - \tau^{\text{mse}}) \xrightarrow{\text{dist}} \mathbf{N}(0, \sigma^2(\tau)).$$

**Question 2.3** (Constructions of conformal confidence sets, 15 points): Suppose we have set-valued mappings  $C_\tau : \mathcal{X} \rightrightarrows \mathcal{Y}$ , meaning that  $C_\tau(x) \subset \mathcal{Y}$ , indexed by  $\tau \in \mathbb{R}_+$ , where

$$C_\tau(x) \subset C_{\tau+\delta}(x)$$

<sup>1</sup>We have  $\mathbb{E}[v^T x (u^T x)^2] = \sum_{i=1}^d v_i \mathbb{E}[x_i (u^T x)^2]$ , and (w.l.o.g. taking  $i = 1$ ) we observe  $\mathbb{E}[x_1 (u^T x)^2] = \sum_{i,j=1}^d u_i u_j \mathbb{E}[x_1 x_i x_j]$ . Then note that  $\mathbb{E}[x_1 x_i x_j] = 0$  for any coordinates  $i, j$  when  $x \sim \mathbf{N}(0, I)$ .

for all  $\delta \geq 0$ , where  $\lim_{\tau \rightarrow \infty} C_\tau(x) = \mathcal{Y}$  (that is, for large enough  $\tau$  the confidence set  $C_\tau(x)$  includes all of  $\mathcal{Y}$ ). Define

$$s(x, y) := \inf \{ \tau \in \mathbb{R} \mid y \in C_\tau(x) \}. \quad (2.1)$$

You are given a sample  $(X_i, Y_i)_{i=1}^n \stackrel{\text{iid}}{\sim} P$  of size  $n$  and define  $S_i = s(X_i, Y_i)$  for each  $i$ , then set

$$\hat{\tau}_n := \text{the } (1 + 1/n)(1 - \alpha) \text{ quantile of } \{S_i\}_{i=1}^n.$$

Let  $\hat{C} = C_{\hat{\tau}_n}$  be the associated confidence set.

- (a) Using the results from class, show that  $\hat{C}$  is a valid  $(1 - \alpha)$  prediction set, that is, on a new example  $(X_{n+1}, Y_{n+1})$  from  $P$ ,

$$P(Y_{n+1} \in \hat{C}(X_{n+1})) \geq 1 - \alpha.$$

We now explore different constructions of such confidence sets. Each of these will leverage an already constructed predictor  $f$  taking inputs in  $\mathcal{X}$ .

- (b) Let  $\ell$  be a loss function and  $\ell(f(x), y)$  be the loss for predicting  $f(x)$  on response  $y$ . Set

$$C_\tau(x) = \{y \in \mathcal{Y} \mid \ell(f(x), y) \leq \tau\}.$$

Give the value  $s(x, y)$  the definition (2.1) yields.

- (c) For binary logistic regression,  $f : \mathcal{X} \rightarrow \mathbb{R}$ ,  $y \in \{\pm 1\}$ , and  $\ell(f(x), y) = \log(1 + e^{-yf(x)})$ . Give the value  $s(x, y)$  the definition (2.1) yields for the confidence set in part (b). For a given threshold  $\tau$ , when is  $C_\tau(x)$  a singleton?
- (d) For  $k$ -class logistic regression,  $f : \mathcal{X} \rightarrow \mathbb{R}^k$ ,  $y \in \{1, \dots, k\}$ , and  $\ell(f(x), y) = \log(1 + \sum_{l=1}^k e^{f_l(x) - f_y(x)})$ . Give the value  $s(x, y)$  the definition (2.1) yields for the confidence set in part (b). For a given threshold  $\tau$ , when is  $C_\tau(x)$  a singleton?
- (e) Let  $\mathcal{Y} = \mathbb{R}$  (so we have real-valued responses as in regression), and let  $l, u : \mathcal{X} \rightarrow \mathbb{R}$  model lower and upper quantiles of  $Y$  given  $X$ , respectively. (That is, we wish to have  $Y \in [l(x), u(x)]$  with a given probability.) Let

$$C_\tau(x) = [l(x) - \tau, u(x) + \tau]$$

where  $C_\tau(x) = \emptyset$  if  $l(x) - \tau > u(x) + \tau$ , i.e., the lower end of the interval is greater than the upper. Give the value  $s(x, y)$  the definition (2.1) yields for this confidence set.

**Question 2.4** (Conformal sets, 20 points): Consider the following heteroskedastic linear model:

$$y = x^T \beta^* + \|x\|_2 \varepsilon, \quad \varepsilon \stackrel{\text{iid}}{\sim} \mathbf{N}(0, 1), \quad x \stackrel{\text{iid}}{\sim} \begin{cases} \mathbf{N}(e_1, I_d) & \text{w.p. } \frac{1}{2} \\ \mathbf{N}(-e_1, 3I_d) & \text{w.p. } \frac{1}{2} \end{cases} \quad (2.2)$$

where  $e_1$  is the first standard basis vector. Given a sample of size  $n$  from this model, let  $X \in \mathbb{R}^{n \times d}$  and  $Y \in \mathbb{R}^n$  be the usual design matrix and responses. Consider fitting the following two models to this data: first, the standard linear regression model

$$\hat{\beta}^{\text{mse}} := \underset{\beta}{\operatorname{argmin}} \|X\beta - Y\|_2^2.$$

This gives predictor  $\hat{f}(x) = x^T \hat{\beta}^{\text{mse}}$ . Second, a quantile model, which attempts to predict the lower and upper  $\alpha$  quantiles of the responses  $y_i$ . To do this, define the feature mapping  $\phi(x) = (x, \|x\|_2) \in \mathbb{R}^{d+1}$ , the quantile loss function

$$\ell_\alpha(t, y) := \alpha(y - t)_+ + (1 - \alpha)(t - y)_+,$$

and then find the  $(d + 1)$ -dimensional vectors  $\hat{\theta}_\alpha$  and  $\hat{\theta}_{1-\alpha}$  solving

$$\hat{\theta}_\alpha = \operatorname{argmin}_\theta \sum_{i=1}^n \ell_\alpha(\theta^T \phi(x_i), y_i) \quad \text{and} \quad \hat{\theta}_{1-\alpha} = \operatorname{argmin}_\theta \sum_{i=1}^n \ell_{1-\alpha}(\theta^T \phi(x_i), y_i).$$

This gives lower and upper predictors  $\hat{l}(x) = \phi(x)^T \hat{\theta}_\alpha$  and  $\hat{u}(x) = \phi(x)^T \hat{\theta}_{1-\alpha}$ .

- (a) What is the population counterpart of  $\hat{\beta}^{\text{mse}}$ ? That is, give  $\beta^* = \operatorname{argmin}_\beta \mathbb{E}[(y - x^T \beta)^2]$ .
- (b) What are the population counterparts of  $\hat{\theta}_\alpha$  and  $\hat{\theta}_{1-\alpha}$ ? That is, give

$$\theta_\alpha^* = \operatorname{argmin}_\theta \mathbb{E}[\ell_\alpha(\theta^T \phi(x), y)] \quad \text{and} \quad \theta_{1-\alpha}^* = \operatorname{argmin}_\theta \mathbb{E}[\ell_{1-\alpha}(\theta^T \phi(x), y)].$$

- (c) Generate three datasets from the model (2.2), each in dimension  $d = 5$  with sample size  $n = 400$ : a training set, a validation set, and a test set, where  $\beta^* \sim \text{Uni}(\mathbb{S}^{d-1})$ . Now, fit the linear predictor  $\hat{f}$  and lower/upper predictors  $\hat{l}, \hat{u}$  on the training data. Consider the confidence sets

$$\hat{C}_\tau^{\text{mse}}(x) := [\hat{f}(x) - \tau, \hat{f}(x) + \tau] \quad \text{and} \quad \hat{C}_\tau^{\text{q}} := [\hat{l}(x) - \tau, \hat{u}(x) + \tau].$$

Using the validation data, use conformal inference to choose  $\tau$  so that  $\mathbb{P}(Y^* \in \hat{C}_\tau(X^*)) \geq 1 - 2\alpha$  for each of these confidence sets, where  $\alpha = .025$ . Repeat this experiment  $T = 100$  times and give the (empirical) coverage you obtain on the test set for each method.

- (d) As in part (c) (with  $T = 100$  experiments), give average the empirical coverage on the following two subsets of the test set:

$$S_{\text{left}} := \{i \mid e_1^T x_i \leq 0\} \quad \text{and} \quad S_{\text{right}} := \{i \mid e_1^T x_i > 0\}.$$

Explain your result in a few words.

**Question 2.5** (A loan data analysis challenge, 40 points): A company in Chile uses crowdsourcing to fund loans to the public, as a means to offer relief from the high bank interest rates. The data in this challenge consists of historical loan records for a sample of 9000 past customers. The variables characterize some aspects of the loan, such as duration, amount, interest rate and many other more technical features of the loans. There are also a number of qualitative variables, such as reason for loan, quality rating of the borrower and others. The response variable  $y$  of interest is `default`: a 0-1 variable indicating whether or not the borrower has defaulted on their loan payments.

The company would like to build a default risk score so that they can target high-risk customers early and perhaps preempt the default event, which ends up costly for all involved. (The fraction of defaults in the entire population is around 7%.) The training data `loan-train.csv` represents a sample of 3000 defaulters, and 6000 non-defaulters, and contains 30 features and the binary outcome `default` (in the first column). The file `loan-testx.csv` consists of a random sample of 10000 other customers from the general pool. For these you are provided only the 30 features.

Your job is to build a *risk score*, that is, a model that estimates the probability of default  $y = 1$ . Feel free to use any of the tools discussed in the lectures of this class (or beyond). Some packages that may be useful include `pytorch`, `xgboost`, and just regular old logistic regression. Describe what you implemented, how you selected your final model. The only thing you need to submit is a text file with 10000 lines; on each line, you should have your predicted risk estimate for each test customer, in the same order as `loan-testx.csv`. Submit this as a `.txt` file on Gradescope for the online portion of HW2.